# Data Needs for LCLS-II

**Amedeo Perazzo**
**SLAC**

U.S. DEPARTMENT OF **ENERGY**
Office of Science

**SLAC** NATIONAL ACCELERATOR LABORATORY
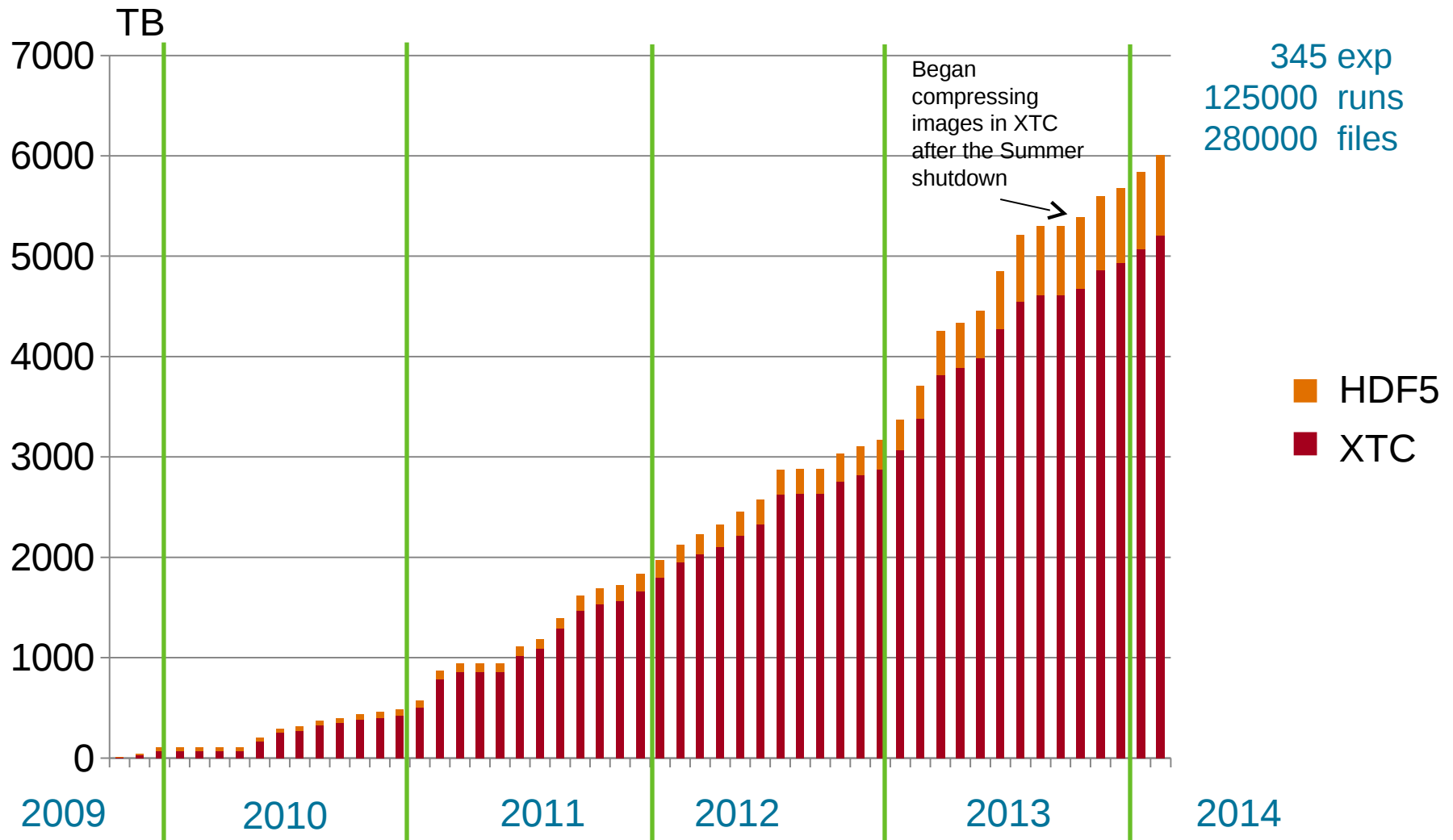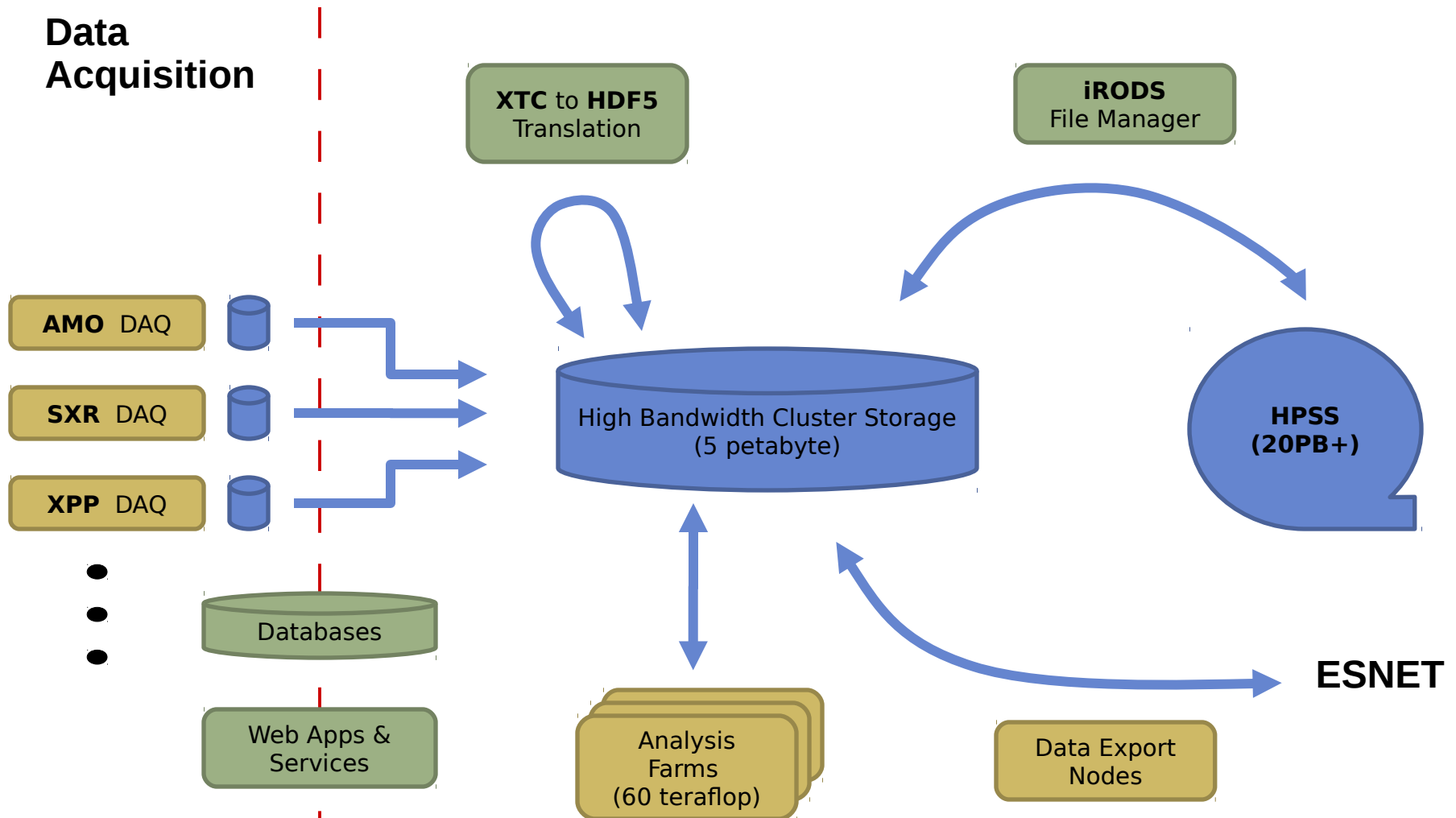
# LCLS Data Throughputs

- **Current LCLS data system can handle fast feedback and offline analysis requirements for most LCLS experiments**

  - DAQ throughput ranges 0.1 – 10 GB/s, typically 1GB/s

    - CSPAD detector: 2 x 2.3 Mpixel @ 120Hz = 1.1 GB/s

- **Predictions for future LCLS data throughput are not obvious**

  - Dictated by project cost, more than physics requirements

  - My guess:

    - One order of magnitude in 4 years time scale

      - 2 x 16Mpixel @ 120Hz (larger CSPAD detectors)

    - Two orders of magnitude in 8 years time scale

      - 100K points @ 100KHz (1D detectors @ LCLS-II data rates)

      - 2 x 4 Mpixel @ 4KHz (ePix detector family)

# LCLS Data Volumes

# LCLS Data Analysis

- **Large variety of tools for analyzing LCLS science data**
  - Real-time, on-the-fly, network based monitoring framework
    - Augmented via modules implemented as shared libraries or shared memory for external framework analysis
  - Fast-feedback, 1-10s delay, disk based analysis
  - Offline analysis: psana (C++/Python), interactive psana, Matlab, CASS, etc
- **Fragmentation of analysis tools partially dictates data infrastructure**
  - Eg. POSIX file systems requirements

# LCLS Data Systems Architecture

**SLAC**

**Data Acquisition**

XTC to **HDF5** Translation

**iRODS** File Manager

**AMO** DAQ

**SXR** DAQ

**XPP** DAQ

High Bandwidth Cluster Storage (5 petabyte)

**HPSS (20PB+)**

Databases

Web Apps & Services

Analysis Farms (60 teraflop)

Data Export Nodes

**ESNET**

# LCLS Data Policies

| Space | Size | Backup | Lifetime | Storage class | Comment |
|---|---|---|---|---|---|
| xtc | Unlimited | Tape archive | 6 months | Short-term | Raw data |
| usr | Unlimited | Tape archive | 6 months | Short-term | Raw data from users' DAQ systems |
| hdf5 | Unlimited | Tape archive | 6 months | Short-term | Data translated to HDF5 |
| scratch | Unlimited | None | 6 months | Short-term | Temporary data |
| xtc/hdf5 | 10TB | n/a | 2 years | Medium-term | Selected XTC and HDF5 runs |
| ftc | 10TB | None | 2 years | Medium-term | Filtered, translated, compressed |
| res | 1TB | Tape | 2 years | Medium-term | Analysis results |
| User home | 20GB | Disk + tape | Indefinite | | User code |
| Tape archive | Unlimited | Two copies | 10 years | Long-term | Raw data |

# LCLS Data Infrastructure

- **DAQ systems dedicated per hutch, user analysis system shared across instruments**
- **Four storage layers**
  - Online cache (flash), fast-feedback (disk), medium term (disk), long term (tape)
  - Medium-term storage currently 5 petabytes
    - Each PB aggregated throughput of 12GB/sec
  - Long-term storage uses tape staging system in the SLAC central computing facilities
    - Can scale up to several petabytes
- **Processing: batch pool and interactive pool**
  - 60Tflop total
  - Most cycles are given out to other SLAC groups because of the bursty nature of LCLS experiments
- **Farms live in the experimental areas with fast (IB QDR) access to the science data files in medium-term storage**

# LCLS Data Management Framework

**SLAC**

- **Data Management system handles all content-opaque operations**

  - Moves data across storage layers (online cache, fast-feedback, offline storage, tape)

  - User accessible through LCLS web-portal (electronic logbook)

  - Handles data policies (security, access, retention)

  - Handles DAQ generated data or data resulted from centralized processing (eg HDF5 translation, compression, filtering)

  - Archive to tape (HPSS) implemented as iRODS service

- **Currently handling 11PB LCLS data, raw and user generated**

  - 5PB on disk, 6PB on tape

# LCLS Data Management Framework
## Interface Examples



Group Management

Translation Manager

Elog Screenshot

File Manager

# Vetoing Events for FEL Experiments Can Be Tricky

- **Very hard to implement effective trigger/veto system**

    - Not strictly a technical issue: the ability to veto events is already implemented in the system

    - Vetoing based on beam parameters not effective (most pulses are good)

    - Hard to get help from users in setting veto parameters which define event quality

        - Users themselves often don't know what these parameters or their thresholds should be

        - Users are usually very suspicious of anything which can filter data on-the-fly

        - Things may get better as algorithms mature

- **Benefit of vetoing events based on the event data is potentially very large for some experiments**

    - Factor 10-100 for some CXI imaging experiments

    - Many experiments, though, have hit rates close to 100%

# LCLS/NERSC Data Pilot

- In 2012 PCDS requested and obtained a NERSC allocation under the "Data Intensive Computing Pilot Program"

- PCDS provided a data-mover script and web-based monitoring to automatically transfer the data for a CXI experiment to NERSC
  - Moved data from SLAC to NERSC at around 700MB/s (ie half of data taking rate)
- PCDS ported LCLS analysis framework to Carver (NERSC farm)

- This exercise showed that partnering with large computer centers like NERSC is part of the solution to LCLS data challenge but can't replace local midscale computing for fast feedback and initial analysis

- Collaborations beyond the data pilot would require 100Gb connection between SLAC and ESNET



Data Transfer from SLAC to NERSC

Feb 28 - Mar 5 Data Transfer

Legend:
- DAQ
- OFFLINE
- NERSC

Shift 5: 29.9 TB
Shift 4: 33.0 TB
Shift 3: 7.1 TB
Shift 2: 28.8 TB
Shift 1: 14.8 TB

Data Progress (TB)

# Offloading LCLS Data Analysis Infrastructure

- **Data centers built towards data intensive systems could help offload the LCLS/SLAC offline computing system**
  - Based on expected data scaling, no modifications to data retention policies, general support for LCLS offline analysis in 2-3 years timescale would require:
    - ~50 PB tape storage, dedicated ~10 PB of disk storage, ~100 teraflop processing farm with an aggregate throughput to the storage above 10 GB/s per PB

- **Key requirements: ability for LCLS users to manage their data through the LCLS tools and workflows, ability to use their SLAC account (or a federated account)**